



Achieving optimal price/performance for your HPC workload

A guide to discovering the best AWS instances
and configurations for your use case

Introduction

1. Key decision points in moving to AWS...
2. Put it all together with AWS ParallelCluster....
3. Control your HPC costs and budget
4. Fine-tune your optimization techniques
5. Achieve optimal price-performance....

Conclusion



INTRODUCTION

High Performance Computing (HPC) drives the frontier of science, finance, engineering, education, chemistry, genomics, modeling, simulation and a host of other compute-intensive industry verticals. The traditional way of running HPC on-premises involves long procurement cycles, high up-front capital investment and, increasingly often, mid-cycle technology refreshes.

The flexibility and heterogeneity of HPC cloud services provide a welcome contrast to the constraints of on-premises HPC. Every HPC configuration is potentially accessible to any given workload in a well-resourced cloud HPC deployment, with vast scalability to spin up as much compute as that workload demands in any given moment.

In 2011, 13% of all sites running HPC leveraged public cloud HPC services. In 2018, the amount rose to 74%.¹ As HPC in the cloud expands into more and more industries and use cases, Amazon Web Services (AWS) is broadening its footprint too. AWS is the primary cloud provider for 58% of the surveyed HPC user organizations running HPC workloads in the cloud, according to Hyperion Research.¹

AWS provides a broad range of HPC capabilities for users' diverse resources, system requirements and data. HPC workloads for which AWS has proven beneficial include:



Tightly-coupled HPC:

Computations involving high numbers of CPUs (e.g. CFD, global climate modeling) readily scale out with AWS to meet expanded demands on the fly. On-premises HPC may not achieve similar flexible performance, due to constraints around available capacity.



Loosely-coupled grid computing:

Some more loosely-coupled workloads (e.g. Monte Carlo calculations, financial risk forecasting and proteomics) may not always require high-performance interconnect or storage, even as the many CPU cores they occupy still test the system's speed limits. On-premises HPC may not be able to adjust to accommodate the workload, leaving system resources underutilized. Spot instances by contrast, opportunistically use spare capacity of only the HPC resources in demand — saving both time and money compared to the on-premises deployment.



High-volume data analytics & interpretation:

Specialized HPC storage like object, block and network file system (NFS) are practical prerequisites for high-volume applications like genomics, high-res image analysis and seismic data processing. Yet any on-premises HPC deployment will have only a single portfolio of storage options — which for any specialized problem could be less than optimal. On the other hand, AWS data storage options are broadly optimizable for each individual data-intensive computing workload.



Accelerated computing:

Some HPC workloads may benefit from the flexibility to activate or increase the mixture of GPU and FPGA accelerators in a compute workload. As with storage options above, an on-premises HPC portfolio will not necessarily “flex” to add additional GPUs or FPGAs as they're needed. AWS resources, by contrast, include a full range of CPU, GPU and FPGA-based instances.



Machine learning and AI:

On-premises HPC administrators today struggle to keep pace with the demands of their growing neural network training and inference workloads. AWS machine learning, deep learning and AI-optimized instances leverage GPU power, specialized FPGA acceleration and high-memory instances for inference — all of which are customizable to fit a user's workload needs.



Data lakes for IoT and analytics:

AWS provides a comprehensive set of services to move, store and analyze your data for data lakes and analytics solutions. AWS supports channel and data stores in buckets. This feature allows users to integrate IoT Analytics data with their existing data lake, to manage the lifecycle of the data according to their bucket policies, and to use the data with a downstream application for further processing or presentation to end-users.

AWS customers benefit from a data center and network architecture built to meet the requirements of the most security-sensitive organizations. This means customers can maintain the highest standards of security without having to manage their own facility. Data privacy is also central to AWS operations. Strong safeguards are in place throughout AWS's infrastructure to help protect customer privacy. (For more, see the callout box: "The High-Performance Security of AWS".)

Intel, a leader in data-centric technologies, has been at the forefront of HPC, delivering products that tackle the workload complexity and challenges of today and tomorrow. AWS, with many different instances built on Intel architecture allows engineers and researchers to develop applications faster and to modernize code with a broad range of optimized software tools, frameworks, and libraries. Together, AWS and Intel provide HPC users ample opportunity to share and collaborate efficiently and securely with team members across the globe.



How To Measure HPC Performance

Technical metrics like performance benchmarks — germane for physical HPC clusters in the data center — can lead away from an optimal cloud HPC mindset. Additional factors must be considered such as ability to scale, time to value, and end-user productivity. In determining the right deployment method, all of the above considerations must be included in measuring HPC performance.

In many cases, HPC in the cloud provides the most efficient and cost-effective solution. The most comprehensive way to determine HPC workload placement is through a thorough analysis of compute costs and human resources cost (researcher/engineering).

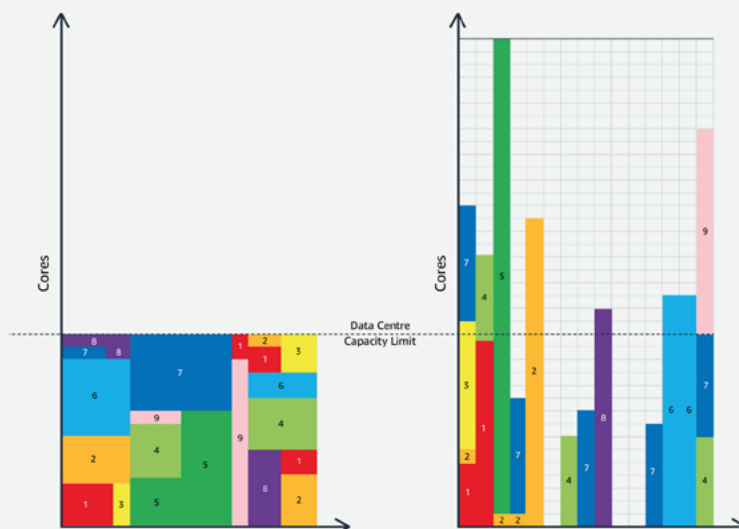


Figure 1. Workloads vs. cores available for an entirely on-premises HPC implementation (Left) and fully-invested HPC cloud solution (Right)

For instance, Team 6 (sapphire blue) in the on-premises HPC scenario spends nine days over two intensive computing periods. The first computing cycle they consume large portions of the HPC cluster's compute — waiting for the job to complete over five days. Then a second job is run which takes four days. Yet the same team only spends two successive days in the HPC cloud scenario. The job clearly requires a high number of cores to complete. But the cloud flexes to accommodate. The HPC cloud scenario frees up seven days of the entire team's expertise, talent and ingenuity that can now be directed to, perhaps, a new design iteration that in-house HPC wouldn't have enabled. Or maybe the team can use the extra time to work on another project altogether.

Meanwhile, Team 8 (purple) occupies what could be a frustrating six days, first skimming off the remaining HPC cluster's compute power and then more than a week later, returning to their task over two further compute days in which they finally have secured enough HPC cycles to finish the job. Contrast that to the single day they run their HPC job across as many cores as their workload demands. The scalability of HPC in the cloud allows Team 8 the same freedom Team 6 experiences.

As this hypothetical but representative example illustrates, the real-world experiences of HPC users in the on-premises HPC implementation are markedly different from the fully-invested HPC cloud solution. After all, what talented engineer, research scientist, R&D developer or other HPC user would prefer to sit in their offices waiting for their job to complete for days on end rather than experience rapid job turnaround? For that matter, what business would want to pay for them to be unproductive?



1. KEY DECISION POINTS IN MOVING TO AWS

A. Moving data to the AWS cloud

The journey begins with moving your data into the cloud. **AWS DataSync** automates moving data between on-premises storage and Amazon S3 or Amazon Elastic File System (Amazon EFS). DataSync automatically handles running your own instances, encryption, managing scripts, network optimization, and data integrity validation. **AWS Direct Connect** (which establishes a dedicated AWS network connection) often enables users to reduce their network costs, increase bandwidth throughput, and provide a more consistent network experience. **AWS Snowball** and **AWS Snowmobile** provide secure transfer of large stores of data (with low impact on network costs and efficient data transfer times) into and out of the AWS Cloud. Customers are encouraged to choose the options that make the most sense for their data and workload.

Services to Get Started with HPC ON AWS

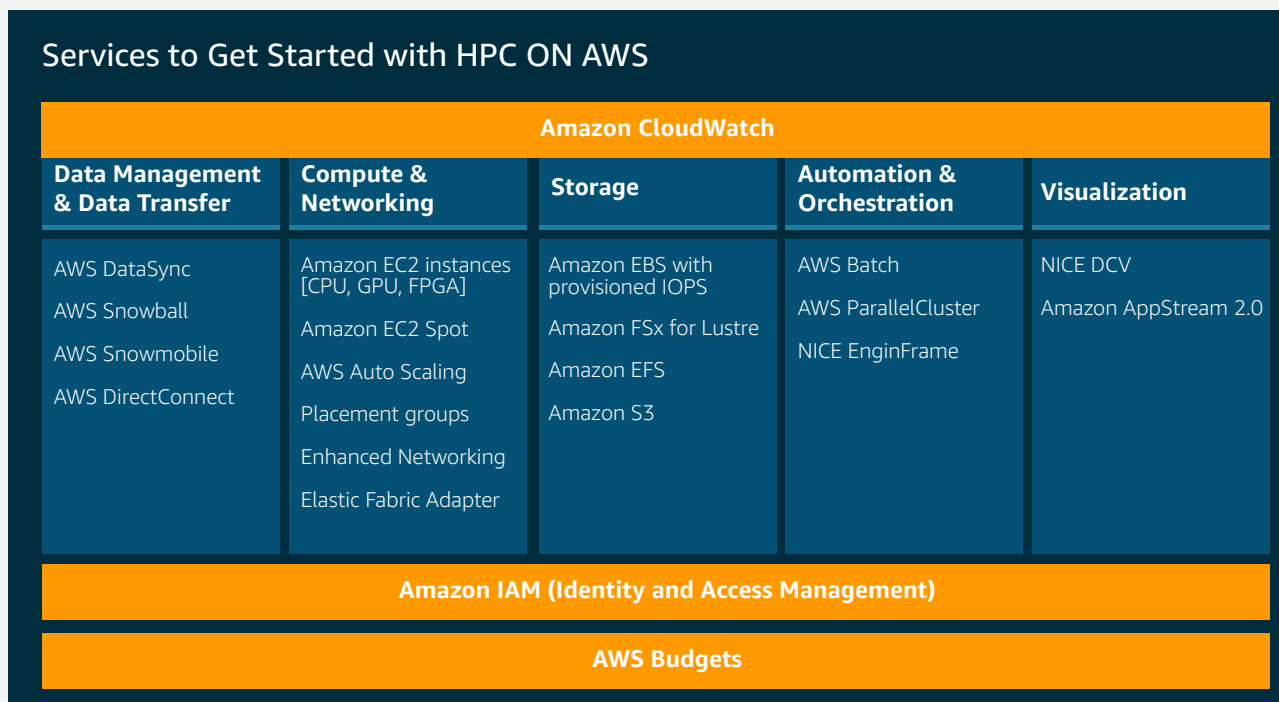


Figure 2. Getting started with HPC on AWS

B. Choosing the right instance type(s) for your compute requirements:

At the core of any HPC solution is flexible compute. With AWS, secure, resizable compute capacity in the cloud comes courtesy of the Amazon Elastic Compute Cloud (Amazon EC2). EC2 instances are optimized for a range of workloads. You can choose the CPU, GPU (or, if relevant, FPGA), memory, storage, and networking capacity that you need to run your applications.

See Figure 2 for the five branches of a standard AWS solution: **Data Management & Data Transfer, Compute & Networking, Storage, Automation & Orchestration and Visualization.**

For most HPC customers, the AWS **C instance** represents the best option for their workloads. “C” here represents CPU or compute-intensive workloads. C5 instances (the current generation) are powered by 3.0 GHz Intel® Xeon® Scalable (a.k.a. Skylake) processors and allow a single core to run up to 3.5 GHz with the Intel® Turbo Boost feature.

Other powerful members of the **C5 family** include the **C5.12xlarge** and **C5.24xlarge** instances run on Intel’s Second Generation Xeon Scalable processors (code-named Cascade Lake) with sustained all-core turbo frequency of 3.6GHz and maximum single core turbo frequency of 3.9GHz.

C5n instances use up to 100 Gbps of network bandwidth with 33% higher memory footprint compared to C5. **C5n** would be the preferred instance when network throughput and packet rate performance are a premium.

Use cases for C5 instances include:

- Computational Fluid Dynamics**
 CFD calculations offer both solid “strong scaling” performance and exceptional “weak scaling” as well. (Strong vs. weak scaling in this context refers to a calculation remaining fixed in grid size or resolution vs. increasing in grid size or resolution as the computer power goes up.) In either case, AWS offers high throughput, security, cost-savings, and availability for the leading CFD software products and ecosystems.
- Tightly-coupled MPI workloads**
 Weather modeling and reservoir simulation are two examples of the kind of tightly-coupled application that can demand high levels of inter-instance communications. AWS’s **C5n instances** offer up to 100 Gbps of network bandwidth, proving particularly useful for simulations, in-memory caches, data lakes and other communication-intensive applications.

The **M family** of AWS instances provide a balance of compute, memory and network resources — and are the option to select for workloads requiring more memory than available on C instance types. M5 has up to 48 physical Xeon cores (96 vCPUs) and 384 GiB of memory (a ratio of 8 GiB memory per physical core). Shared memory applications, for instance, or other tightly coupled workloads may find particular use for M5’s large number of cores. This enables, for instance, manufacturing organizations to run their Finite Element Analysis (FEA) and simulation jobs to completion more quickly.

AWS instances in the **R family** represent memory-intensive jobs. R5 instances have 1:8 vCPU to memory ratio, with the largest offering as much as 768 GiB per instance. These instances feature Intel’s Xeon Platinum 8000 series (Skylake-SP) with a sustained all-core Turbo CPU clock speed of up to 3.1 GHz.

Instance Types	M5	M5d	R5	R5d	C4	C5	C5d	C5n	Z1d	P3	P3dn	G3	F1
Example Use Case	FEA Implicit		CFD, FEA Explicit				EDA, CFD	ML/AI CUDA		Remote Visualization	Genomics, Finance		
Max CPU (GHz)	3.1		3.1		2.9	3.9	3.5		4.0	2.7	2.7	2.7	2.7
Max RAM (GB)	384		768		60	192	144	192	384	488	768	488	976
Max vCPUs	96		96		36	96	72	72	48	64	96	64	64
Max Cores (*)	48		48		18	48	36	36	24	32	48	32	32
RAM/vCPUs	4		8		1.6	2		2.6	8	7.6	8	7.6	15.25
RAM Cores (*)	8		16		3.3	4		5.3	16	15.2	16	15.2	30.5
Max NVME SSD (TB)	NA	1.8	NA	3.6	NA	NA	1.8	NA	1.8	NA	1.8	NA	3.7
Max Network Bandwidth (Gbps)	25	25	25	25	10	25		100	25	25	100	25	25
Accelerated Computing	—									Up to 8 Nvidia Volta V100		Up to 4 Nvidia Tesla M60	Up to 8 Xilinx FPGAs

Figure 3. Comparison of AWS instance types and example use cases for those instances. Please note that because EC2 instances support multithreading, multiple threads run concurrently on a single CPU core. Each thread is represented as a virtual CPU (vCPU) on the instance. (Each vCPU is a thread of a CPU core, except for T2 instances.)

Z instances provide high single-thread performance courtesy of a custom Intel Xeon Scalable processor (exclusive to AWS) that achieves up to 4 GHz sustained, all-turbo performance, enhanced networking and up to 25 GB throughput.

Z1d instances work well with memory and compute-intensive applications that require high-frequency processing. For instance, Z1d users leverage the sustained core performance for applications with license restrictions that require single-threaded applications and those whose software is licensed by the core.

For all these instances, AWS and Intel also enable faster HPC app development and code modernization with a broad range of optimized software tools, frameworks, and libraries on AWS.

F1 instances provide HPC users access to high speed field programmable gate arrays (FPGAs, a kind of specialty integrated compute engine used, for instance, in specialized workloads in finance and genomics). One genomics application uses F1 instances for HPC genome sequencing and drug targeting. This customer reports being able to analyze 100,000 exomes in 100 hours using AWS — providing among the fastest and most efficient genomics pipelines in the industry.

For HPC users seeking GPU accelerators, AWS offers both the **P and G family** of instances for graphic intensive or CUDA applications.

P3 instances accelerate HPC workloads via powerful, industry-leading GPUs. P3dn instances are optimized for distributed machine learning training with Amazon Machine Images (AMI) — including most of the popular machine learning frameworks as well as Amazon Sagemaker.

Industries that have used Amazon EC2 P3 instances to boost their core HPC capabilities include:

- **Data Storage - Hard Disk Design**
Companies use P3 instances to run material science, thermodynamic, magnetic and data transfer simulations at high speeds.
- **Pharmaceutical**
Companies leverage P3 instances to extend the scale of discovery and accelerate their drug development pipeline — increasing the number of simulations and lines of investigation their researchers can pursue.
- **Seismic processing**
Teams in oil & gas companies can demand unique CPU and GPU configurations — with a scale and elasticity to support spiky optimization workflows, like automated history-matching. P3 instances can help engineers iterate and fine-tune their models faster, thereby accelerating reservoir simulations.

Industries that have leveraged G3 instances include:

- **Energy exploration and production**
In which G3 instances power increasingly complex exploration and production (E&P) as well as overall reservoir management.
- **Media and entertainment**
Whose graphics-intensive workloads include post production, video playout/broadcast, video encoding transcoding and AR/VR applications.

See Figure 3 for the full specs on each of AWS's HPC instances.

C. Selecting the storage your workload will need

For HPC workloads requiring block storage, which can include in-memory databases and genomics workloads, Amazon Elastic Block Store (EBS) provides persistent block storage for read-often temporary working storage and high-IOPS tasks. EBS optimization —enabled by default on C5 and M5 instances — minimizes contention for the network, allowing better and more consistent application performance.

EBS provides a number of options to enable storage-optimized performance that minimizes cost. EBS General Purpose SSD (gp2) volumes work best with workloads without a continuous I/O demand. More I/O intensive applications are better situated with Provisioned IOPS SSD (io1).

For workloads with high I/O needs, AWS has developed a HPC-optimized parallel file system called Amazon FSx for Lustre. FSx for Lustre offers massively scalable performance with consistent sub-millisecond latencies. Each TB provides 200 MB/second and scales to hundreds of GB/s and millions of IOPS.

AWS also provides HPC storage volumes for object storage (S3), file system storage (EFS) and archival storage (Amazon Glacier) for low-cost data archive and backup.

As an example, financial services industry applications can involve as much as trillions of individual records to monitor and run database queries against. One industry regulatory body (FINRA - Financial Industry Regulatory Authority) uses S3 storage as part of a market regulation portfolio that includes EC2 instances running Hadoop and Apache HBase. The compute instances provided sub-second to minutes wait times for queries (compared to minutes to hours for their previous on-premises solution). The S3 storage provided flexibility of switching EC2 instances that are optimized for each query or workload. This way the user did not have to pay multiple times for storing enormous data files on separate EC2 instances — as well as not having to transfer large data volumes to new instances as they're spun up.²

D. Choosing the networking best suited to your workload

When it comes to choosing networking options for your workloads, AWS provides a number of options:

Elastic Network Adapter (ENA): Enhanced Networking, provides higher I/O performance and lower CPU utilization when supporting network traffic. Modern instance types can benefit from the Elastic Network Adapter (ENA). ENA is a custom network interface optimized to deliver high throughput and packet per second (PPS) performance. ENA supports network speeds of up to 25 Gbps. The latest Amazon Linux HVM AMIs have the module required for enhanced networking with ENA and have the required attributes set. Therefore, if you launch an instance with the latest Amazon Linux HVM AMI on a supported instance type, enhanced networking is already enabled for your instance. If you are not using the standard Amazon Linux AMI, please follow [this guide](#) to configure your environment.

Elastic Fabric Adapter (EFA): EFA is a network interface for Amazon EC2 instances that enables customers to run HPC applications requiring high levels of inter-instance communications, like computational fluid dynamics, weather modeling, and reservoir simulation, at scale on AWS. It uses a custom-built operating system bypass technique to enhance the performance of inter-instance communications, which is critical to scaling HPC applications. With EFA, HPC applications using popular HPC technologies like Message Passing Interface (MPI) can scale to thousands of CPU cores. EFA supports industry-standard libfabric APIs, so applications that use a supported MPI library can be migrated to AWS with little or no modification. EFA is available as an optional EC2 networking feature that you can enable on supported EC2 instances like C5n.18xlarge, P3dn.24xlarge and i3en.24xl. [Here](#) you can find how to configure EFA.

E. Deciding on your automation and orchestration preferences

Running AWS instances often involves schedulers and batch management systems. Users can bring their own scheduler to AWS or develop their own solutions (e.g. using SQS and Cloudwatch to monitor and signal autoscaling, to add resources on demand or to use the spot market).

On the other hand, external schedulers may not always be necessary. **AWS Batch** is a managed job scheduler that dynamically provisions compute resources based on the volume and requirements of the submitted jobs. It will plan, schedule and execute batch jobs across cloud HPC resources.

AWS provides other tools for automation and orchestration of HPC jobs, including ParallelCluster, which enables the launch of virtual HPC clusters in minutes (see Section 2 below). CloudFormation templates can also be built to automate the entire HPC workload.

F. Incorporating the right visualization tools for your HPC output

Visualization of results is an important aspect of a HPC workflow. Remote visualization helps accelerate turnaround times for engineers as well as team managers.

NICE Desktop Cloud Visualization, for instance, is a remote visualization technology that enables users to securely connect to graphic-intensive 3D applications hosted on a remote high-performance server. NICE DCV makes a server's high-performance graphics processing capabilities available to multiple remote users by creating secure client sessions.

In a typical NICE DCV scenario, a graphic-intensive application, such as a 3D modeling or computer-aided design application, is hosted on a high-performance server that provides a high-end GPU, fast I/O capabilities, and large amounts of memory. The NICE DCV client remotely connects to the session and uses the application hosted on the server. The NICE DCV server software compresses the visual output of the hosted application and streams it back to the user as an encrypted pixel stream.

In addition to NICE DCV, there are other remote visualization services available in AWS:

- **Graphic EC2 Instances**
Amazon EC2 G3 instances are the latest generation of Amazon EC2 GPU graphics instances that deliver a powerful combination of CPU, host memory, and GPU capacity. G3 instances are ideal for graphics-intensive applications such as 3D visualizations, mid to high end virtual workstations, virtual application software, 3D rendering, application streaming, video encoding, gaming, and other server-side graphics workloads. You can install DCV on a G3 instance and start immediately to stream the remote desktop to your local clients.
- **Elastic Graphics**
Amazon Elastic Graphics it's an EC2 feature that allows you to easily attach low-cost graphics acceleration to a wide range of EC2 instances over the network.

Simply choose an instance with the right amount of compute, memory, and storage for your application, and then use Elastic Graphics to add graphics acceleration required by your application for a fraction of the cost of standalone GPU instances such as G3 instances. NICE DCV will automatically benefit of the OpenGL acceleration when installed on an instance with Elastic Graphic.

Amazon AppStream 2.0 is a fully managed application streaming service. You centrally manage your desktop applications on AppStream 2.0 and securely deliver them to any computer. Each scientist or engineer can have a fluid and responsive experience, including GPU-intensive 3D design, because your applications run on virtual machines (VMs) optimized for specific use cases and each streaming session automatically adjusts to network conditions.

2. PUTTING IT ALL TOGETHER WITH AWS PARALLELCLUSTER

AWS ParallelCluster is an AWS-supported open source cluster management tool that helps you to deploy and manage High Performance Computing (HPC) clusters in the AWS cloud.

Built as an enhancement to and replacement for the popular open source CfnCluster project, AWS ParallelCluster enables customers to quickly build a HPC cluster on AWS. It automatically sets up the required compute resources and shared file systems and offers a variety of batch scheduler options, including AWS Batch, Sun Grid Engine (SGE), Torque and Slurm.

AWS ParallelCluster reduces the operational overhead of cluster management and simplifies running HPC workloads on AWS. AWS ParallelCluster facilitates both quick-start proof of concepts (POCs) and production deployments. AWS ParallelCluster is available at no additional charge, and you pay only for the AWS resources needed to run your applications. AWS ParallelCluster is distributed as a Python package and is installed using *pip*.

From here six simple steps remain to access the power of HPC on AWS:

- [Install AWS ParallelCluster](#) using CLI prompts
- [Configure AWS credentials](#), ensuring the user has the proper permissions and storage settings
- [Configure and launch AWS ParallelCluster](#) via a config file that can be customized to toggle features like EC2 placement groups
- [Submit and run a simple parallel MPI job](#) by creating an executable, creating the job submittal file and launching the job
- [Create an Amazon EBS volume snapshot for cluster reusability](#), enabling repeatability and the deployment of the same pre-configured software on future clusters
- [Delete and clean up the cluster](#), making the highly portable AWS HPC experience as easy to tear down as it was to set up.

For more technical, step-by-step instructions and troubleshooting guides, see online resources for AWS ParallelCluster including [“Deploying an Elastic HPC Cluster”](#) (on which the above is based).

3. CONTROL YOUR HPC BUDGET

AWS provides a variety of subscription and pricing options to customize a user's experience and ensure they only pay for the AWS services they need. The AWS Monthly Calculator (<https://calculator.s3.amazonaws.com/index.html>) provides an estimate of usage charges for AWS services based on input information such as data transfer rates and elastic IPs.

AWS HPC users enjoy three flexible payment models, providing maximum output with minimum outlay:

a) Urgent and high-priority workloads

(On Demand Instances): By-the-hour pricing with no long-term commitments or upfront payments, only paying for the instances you use. Payment increments can be by the hour or by the second (with minimum purchase of 60 seconds) — and no long-term commitments.

b) Flexible start and end-time workloads

(Spot Instances): Buy unused EC2 capacity (on average 50%-90% lower than on-demand pricing). Supply and demand determines the price. Spot prices are generally lower than On Demand and Reserved Instances (below). A flipside to the lower price, of course, is the proviso that spot instance jobs may sometimes need to be interrupted. To ensure you have the cost-optimum mix of run-time and price savings, AWS provides the Spot Instance Advisor (<https://aws.amazon.com/ec2/spot/instance-advisor/>).

c) Steady-state usage workloads

(Reserved Instances) provides the confidence AWS HPC users need for the instances they require at the time when they need them.

Reserved Instances are available in three options: All up-front (AURI), partial up-front (PURI) or no upfront payments (NURI). As a general rule, the larger the up-front payment, the greater the discount. So, for example, an instance that costs \$955 per year on demand would only command \$650 per year NURI, \$554 per year PURI and \$545 per year AURI.³

HPC on AWS provides the resources and mechanisms for exerting full control over your HPC workloads — and your HPC budget. As noted in the accompanying sidebar (“How to Measure HPC Performance”), it's important also to consider broader questions of human resources when assessing the overall price-performance of HPC within an organization.

In other words, an on-premises HPC cluster may turn around strong performance metrics and benchmarks in the abstract. But if an organization's researchers, engineers, designers, coders and other HPC users are still spending hours or days waiting for their results using on-premises HPC infrastructure, overall HPC resources may still need adjusting.

HPC in the cloud of course allows users to choose any type of instance and storage configuration your workload might demand. Thousands of cores — even a million cores — can also be spun up in short time frames to handle the scale of variable workloads.

Most important, as readers weigh the various instances and AWS configurations discussed in this white paper, they may have a list of multiple instances that might work well with their suite of applications, data and workloads. This is perfectly normal and in fact is almost to be expected.

Fortunately, AWS provides the ultimate HPC platform for experimenting and testing out any number of configurations without ever committing to future compute configurations or physical infrastructure.

In one word: Experiment. HPC clusters have never been easier to assemble and disassemble, virtually and with just a few command-line commands. Every AWS HPC client should be leveraging AWS's flexibility and scalability wherever and whenever they may need.

4. FINE-TUNE YOUR OPTIMIZATION TECHNIQUES

To optimize and scale your HPC instance, a number of best practices can improve workload performance:

- **Use Intel hardware features.** When compiling applications code, Intel® Advanced Vector Extensions (AVX/AVX2) can be very useful. C5, M5 and R5 instances also provide support for the new Intel Advanced Vector Extensions 512 (AVX-512) instruction set. One example compiler flag that enables AVX-512 is “-xCOMMON-AVX512” — although flags can vary depending on the compiler.
- **Use Cluster Placement Groups.** Tightly-coupled HPC applications often require a low-latency network connection between compute nodes for best performance. On AWS, this is achieved by launching compute nodes directly into a Cluster Placement Group. As the name suggests, this feature clusters the compute nodes close to each other to achieve consistent latency. For maximum effect, launch all compute nodes into a placement group, all at once. All instance types that support enhanced networking can be launched within a Cluster Placement Group. Placement Groups allow for reliably low latency between instances and will help your tightly-coupled application to be elastically scalable as desired. For more information, see the related [AWS documentation here](#).
- **Disable Hyper-Threading.** Amazon EC2 instances support Intel® Hyper-Threading Technology, which enables multiple threads to run concurrently on a single Intel Xeon CPU core. Each thread is represented as a virtual CPU (vCPU) on the instance. An instance has a default number of CPU cores, which varies according to instance type. Except for T2 instances, each vCPU in AWS is a hyperthread of an Intel Xeon CPU core. Most HPC platforms have Intel hyperthreading disabled by default. Unless an application has been thoroughly tested in the hyper-threaded (HT) environment, it's recommend to disable hyper-threading.

For additional information, see this [AWS documentation page](#) and post about [disabling HT on Linux](#) and on [Windows](#).

- **Use real-world data for your tests.** The foremost method to check an application's performance on AWS is to run a meaningful demonstration of the application itself using the same input file you use on your environment. An inadvertently small or large demonstration case, one that does not match expected compute, memory, I/O data transfer or network traffic loads, will not provide a meaningful example of how an application runs on AWS.
- **Experiment. Leverage the power of cloud HPC.** Unlike on-premises HPC, AWS HPC enables users to take multiple combinations of instances out for “test drives” on their workloads. The ability to switch and mix instances on the fly, responding to the demands of each workload (and scaling to many cores in the moment as well) is one of the benefits of cloud HPC for which there is simply no comparable experience in a strictly on-premises HPC procurement.

5. ACHIEVE OPTIMAL HPC PRICE/PERFORMANCE

HPC on AWS leverages the power of cloud computing at the high-performance scale to achieve optimal HPC price/performance. The key for most HPC workloads is for the user to try different instance types and compare their performance with their costs.

With AWS you can right size your services to meet exactly the capacity requirements you need without having to overprovision or compromise capacity. It's easy to choose services that meet your existing workload needs, and as your demands change, you can quickly shift to the services option that meets your new requirements. You can also run multiple service options concurrently, helping you reduce costs and still maintain optimal performance.

One of the first questions that veteran AWS HPC users ask — and that may escape the notice of newcomers to AWS — is how much of any given workload is “spot-friendly”?

Spot Instances can save a user up to 90% of costs by bidding on spare EC2 instances. This way users can realize cost savings for batch processing applications while still maintaining high availability.

When you request a Spot instance, Spot will default the maximum price you are willing to pay per Spot instance-hour as the On-Demand price. You can also exercise additional control over your Spot instance budget by specifying the maximum price you are willing to pay per instance-hour in your request. You will continue to pay the Spot price that's in effect for the time period your instances are running. If Spot price rises above your maximum price, your instance will be automatically terminated, stopped or hibernated.

You can use **Amazon CloudWatch** to collect and track metrics, monitor log files, set alarms, and automatically react to changes in your AWS resources. With Trusted Advisor you can further provision your resources following best practices to improve system performance and reliability, increase security, and look for opportunities to save money.

AWS Cost Explorer gives you additional ability to analyze your costs and usage. Using a set of default reports, you can quickly get started with identifying your underlying cost drivers and usage trends. Cost Explorer enables the user to filter and group data according to resource tags. Tagging makes it easier to map resources and workloads to the appropriate cost center.

One final important price/performance question concerns scale. Namely, does a user's code scale to take advantage of the virtually unlimited capacity AWS can offer? “Lifting” and “shifting” to the cloud, by itself, only provides so much advantage to the HPC customer. What's needed, instead, is a new mindset in which, say, a million cores could actually be mustered to compute on your workload and achieve results and insight on timescales that might previously have been unimaginable.

HPC on AWS provides a cloud environment that embraces a culture of innovation and experimentation and pushes the boundaries of what's possible. AWS customers are taking the cloud conversation beyond simple cost reduction to how businesses can create advantages in an increasingly disruptive economy.



The High-Performance Security of AWS

The AWS infrastructure has been designed to be one of the most flexible and secure cloud computing environments available today. It provides an extremely scalable, highly reliable platform that enables customers to deploy applications and data quickly and securely.

AWS's infrastructure is built and managed not only according to security best practices and standards, but also with the unique needs of the cloud in mind. AWS uses redundant and layered controls, continuous validation and testing, and a substantial amount of automation to ensure that the underlying infrastructure is monitored and protected 24/7. AWS ensures that these controls are replicated in every new data center or service.

AWS's Web Application Firewall (WAF) protects web applications from common exploits via Amazon virtual private cloud. Users can create private networks and restrict access to instances and applications. Web security via WAF is available along the development chain, from the developer to the engineer deploying the software to the security consultants performing security audits.

Transport layer security (TLS) provides industry-standard cryptographic protocols protecting data in transit across all AWS services.

AWS storage and database services encrypt data at rest across the gamut of AWS offerings including EBS, S3, Glacier, Oracle RDS, SQL Server RDS, and Redshift.

AWS Direct Connect allows users to simply establish private or dedicated network connections from premises to AWS. With multiple virtual interfaces, Direct Connect establishes private connectivity to multiple VPCs while maintaining network isolation.

Identity and Access Management (IAM) service provides secure, managed access to AWS resources and services. IAM enables fine-grained access to resources, multi-factor authentication for highly privileged users and access control for mobile applications.

AWS monitoring provided by Amazon CloudWatch — monitors resources as well as provides system-wide visibility into resource utilization, application performance and operational health. CloudWatch issues alert notifications when specific events occur, or thresholds are exceeded.

Machine learning provides security at scale with Amazon Macie and Amazon GuardDuty. Intelligent security firewalls and threat detection measures continuously monitor for malicious or unauthorized behavior, protecting your AWS accounts and workloads. Macie provides a security services that leverages machine learning to automatically discover, classify and protect sensitive data in AWS.

Amazon Quick Start architecture helps support regulatory frameworks such as HIPAA. Quick Start automatically configures AWS resources to generate a complete, secure and compliant environment in minutes.



CONCLUSION

High Performance Computing is more than achieving higher performance levels in compute workloads. It's about delivery of tangible business value in the most cost-effective manner. No longer should HPC deployment decisions be made solely on judging performance benchmarks like LINPACK but rather the speed of iteration for your team. HPC in the cloud delivers the results that matter.

AWS and Intel help you deliver results faster by enabling access to a large catalog of cloud based HPC services to transform your business processes and accelerate the pace of your innovation.

As HPC in the cloud increases its global profile AWS remains the standard bearer. Migrating HPC workloads to AWS represents a continued pathway to increased productivity, access to the latest Intel technologies, scalable results and innovation without constraints.

AWS has a nearly endless configurability of instances and solutions to cater to every HPC user's individual workload. AWS's broad scalability and flexibility provides virtually unlimited infrastructure and fast networking. Whether an organization is switching their HPC over to all-cloud or to a Hybrid HPC middle ground, AWS provides a broad range of scalable, flexible infrastructure services that you can select to match your workloads and tasks. Create an AWS account today and start AWS ParallelCluster to run benchmarks on your workload — and discover how rapidly HPC on AWS will change the game for your business.

Learn more about running your HPC workloads on AWS at <http://aws.amazon.com/hpc>

ACKNOWLEDGMENTS

This white paper gratefully acknowledges previous work by Brendan Bouffler, Francesco Ruffino and Bala Thekkedath.

REFERENCES

¹ Steve Conway, Alex Norton, Bob Sorensen, Earl Joseph, "Cloud Computing for HPC Comes of Age," Hyperion Research, March 2019.

² AWS, Low-Latency Access on Trillions of Records: FINRA's Architecture Using Apache HBase on Amazon EMR with Amazon S3, November 2016. <https://aws.amazon.com/blogs/big-data/low-latency-access-on-trillions-of-records-finras-architecture-using-apache-hbase-on-amazon-emr-with-amazon-s3/>

³ AWS, AWS Pricing, July 2019. https://aws.amazon.com/pricing/?nc2=h_ql_pr